

# Traffic Intensity Model

Google Maps images with traffic layer/colouring: green, orange, red, dark red  
 % of coloured pixels in annular sectors -> traffic intensity; physical density “field”

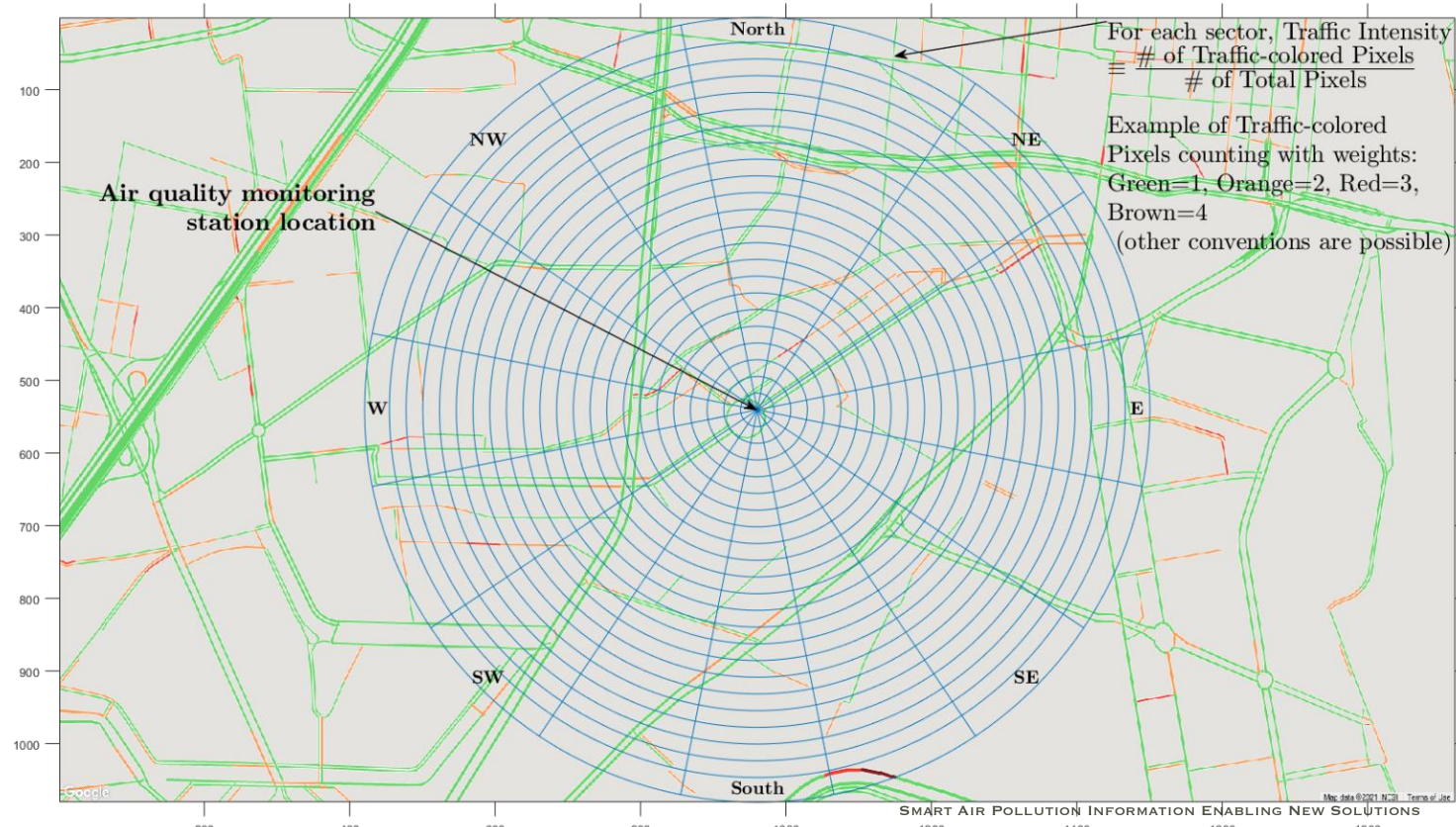
Dimension reduction  
 from HD image 1920x1080:

(4 colours)x(16 angles)x(118 rings)

Or aggregating the angles:

(4 colours)x(118 rings)

Only 23 rings in this plot  
 118 rings shown in the next page



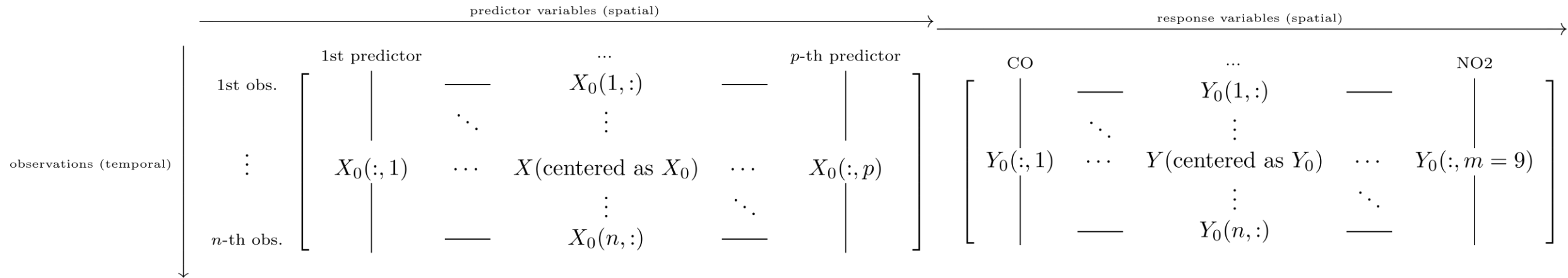
# Traffic Intensity Model



Width of 118  
concentric  
rings: 10m

# Regression Modeling: data matrices

Spatio-Temporal Data Matrices  $X$  (traffic predictors; centered as  $X_0$  with 0 mean & rescaled with Var=1) and  $Y$  (pollutants responses; centered & rescaled as  $Y_0$ ):



where  $p = (4 \text{ colors}) \times (118 \text{ rings})$  (or  $(4 \text{ colors}) \times (16 \text{ angles}) \times (118 \text{ rings})$ ) traffic predictor variables,  $m = 9$  pollutant response variables;

$n$  observations, depending on station/sensor (each has different # of missing/null readings)

# Regression Modeling: PLSR

Purpose of our modeling is two-fold:

- (1) (interpretations) get traffic activities  $\xrightarrow{\text{mapping}}$  pollutant concentrations, and reveal insight on their detailed relations, and
- (2) (predictions) predict pollutant concentrations based on traffic activities.

Most black-box machine learning techniques are good at (2) only;

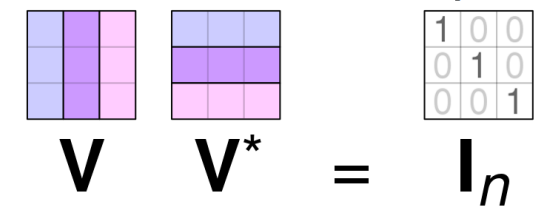
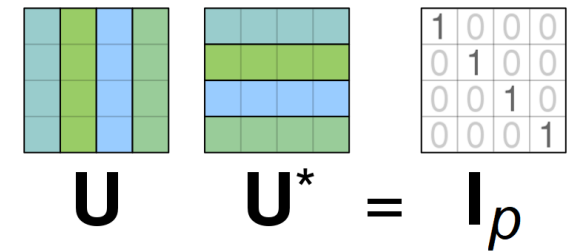
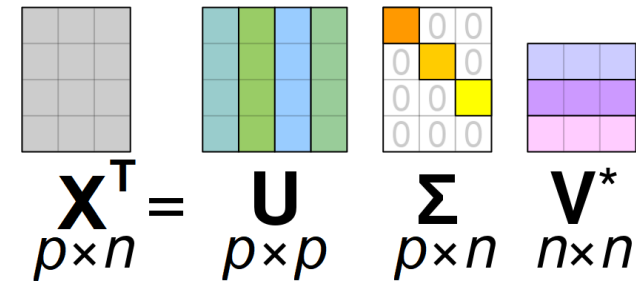
Partial least squares (PLS) regression formulation:

$$\begin{cases} X_0 = X_S X_L^T + X_{\text{residuals}} \\ Y_0 = Y_S Y_L^T + Y_{\text{residuals}} \end{cases}$$



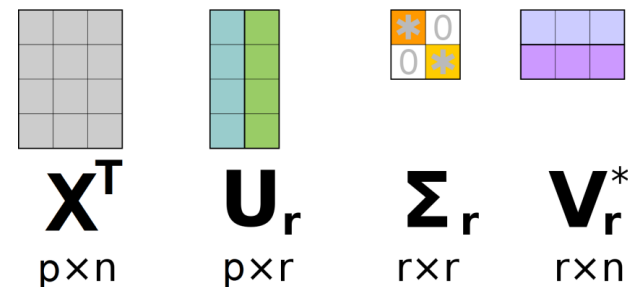
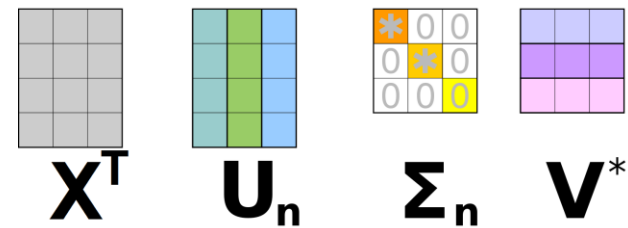
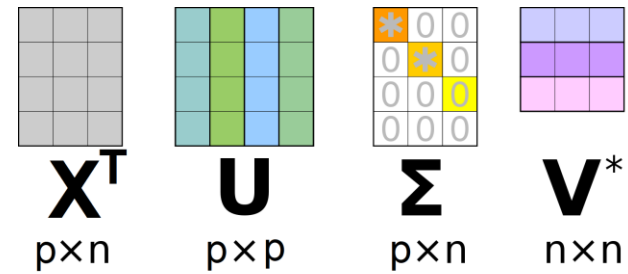
# Preliminary: Singular Value Decomposition (SVD)

- Singular value decomposition:  $X^T = U\Sigma V^*$
- For real matrices,  $U^* = U^T$ ,  $V^* = V^T$ , and
- $U$  &  $V$  are both orthogonal matrices:
- $p$  (spatial; variables)  $>$   $n$  (temporal; observations), reduced



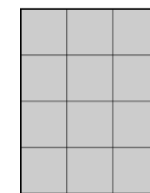
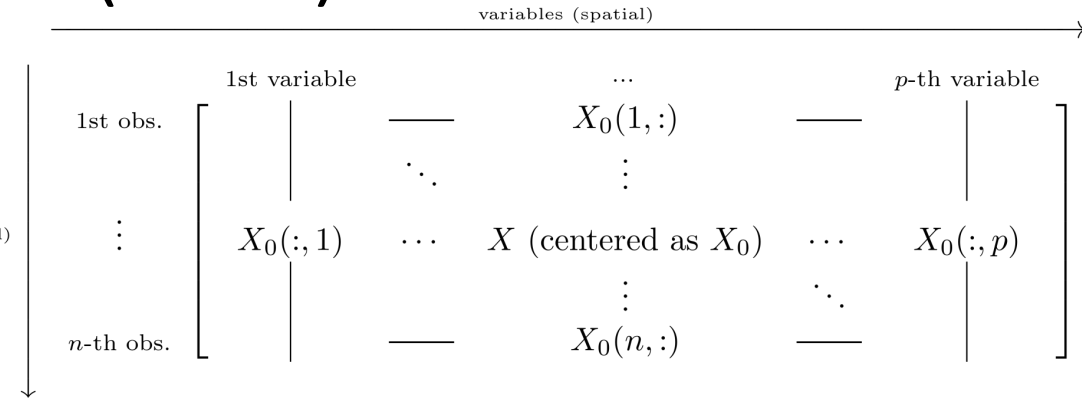
SVD:

- Can reduce further:
- Diagonal line of  $\Sigma$ : non-increasing, truncate the diagonal  $\Sigma$  to 1<sup>st</sup>  $r$  values, and the columns of  $U$  and  $V$
- Columns of  $U$  are still orthonormal basis; same for  $V$

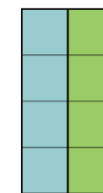


# Principal Component Analysis (PCA)

- Recall the spatio-temporal data matrix:
- Variables (columns) are often correlated
- Covariance matrix of variables  $\propto X^T X$
- Using SVD:  $X^T X = U \Sigma V^T V \Sigma^T U^T = U \Sigma^2 U^T$
- $\Sigma^2 \propto$  eigenvalues of covariance matrix of variables (columns in  $X$ )
- Columns in  $U$  are the corresponding eigenvectors, called Principal Components
- $X^T = U \Sigma V^T$ , so  $U^T X^T = \Sigma V^T$ , i.e., projecting each observation/data point to columns of  $U$ , and get  $\Sigma V^T$  of  $r \times n$  (truncated to  $r$  components)
- Rows of  $\Sigma V^T = X_S^T$  called "scores", capturing major variances up to  $r$  elements of  $\Sigma^2$
- $\|X_S(:, i)\|^2 = \Sigma^2(i, i)$



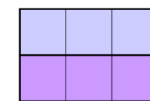
$X^T$   
p x n



$U_r$   
p x r



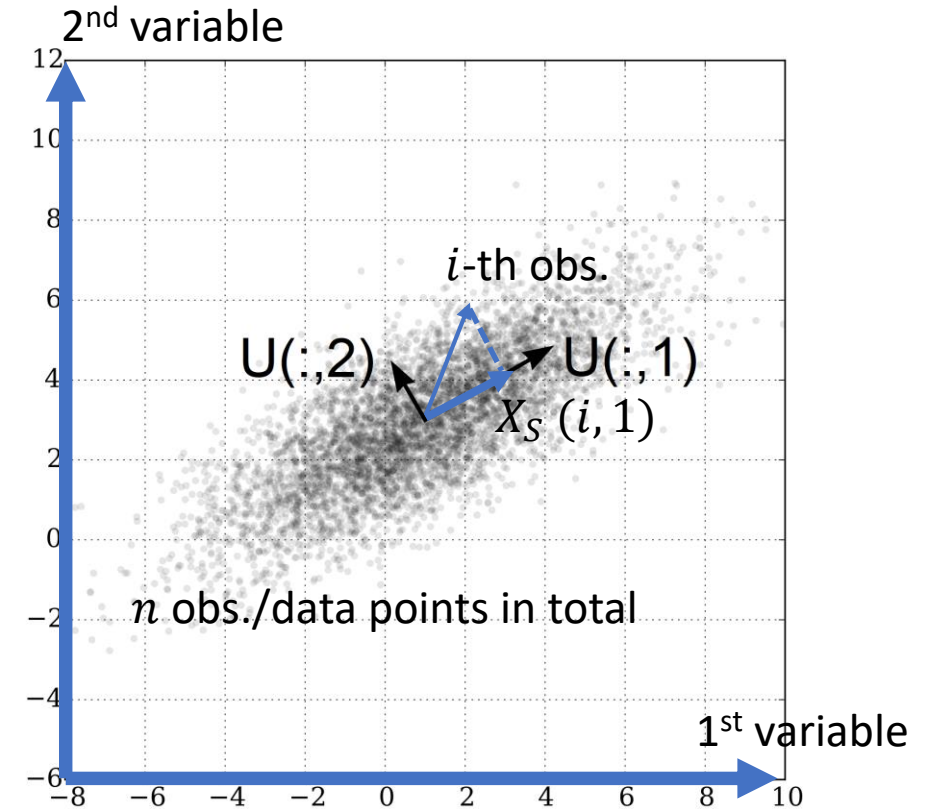
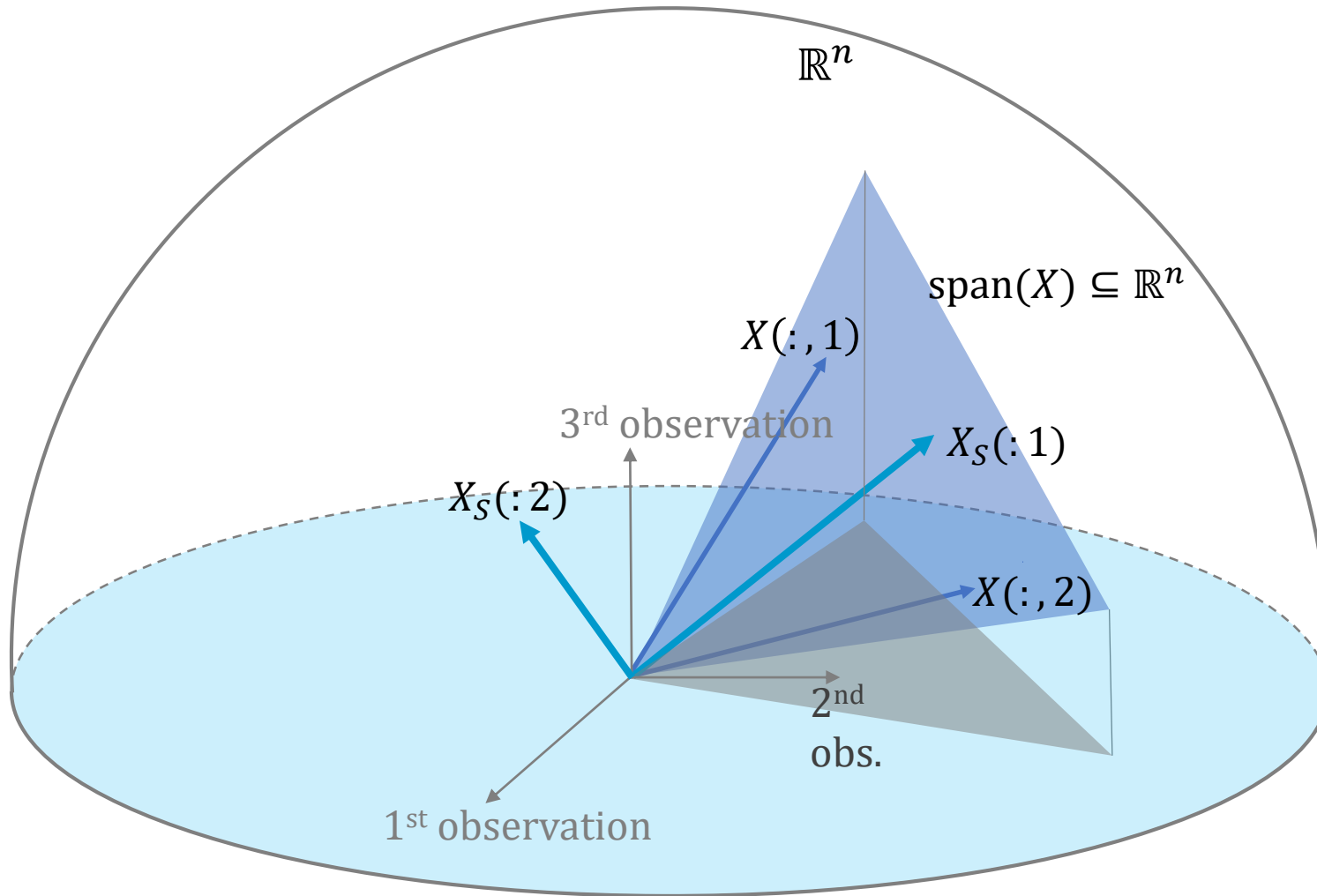
$\Sigma_r$   
r x r



$V_r^*$   
r x n

# Principal Component Analysis (PCA)

- Geometric interpretation:

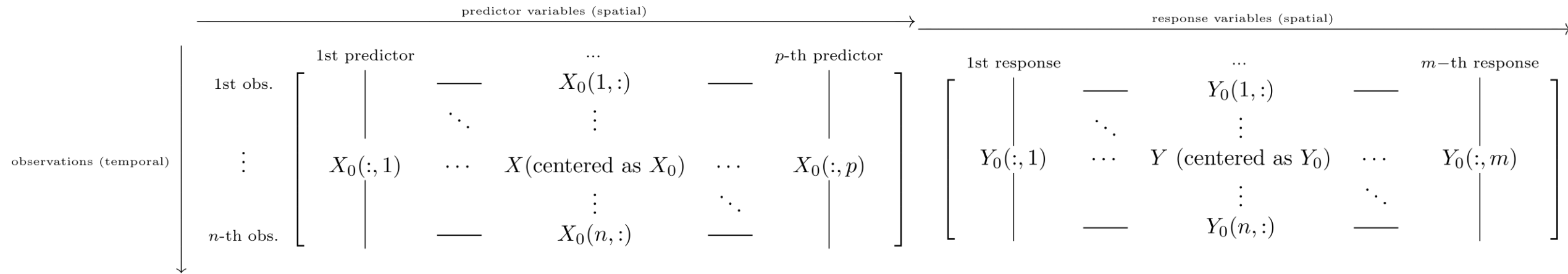


$\|X_S(:, i)\|^2 = \Sigma^2(i, i)$   
Non-increasing along diagonal of  $\Sigma$   
Truncate up to  $X_S(:, r)$  to reduce dimension



# Partial Least Squares (PLS)

- PCA only considers 1 spatio-temporal data matrix
- In regression analysis, we have both predictor variables (predictors, independent variables;  $x$ ) and response variables (responses, dependent variables;  $y$ ), organized in  $X$  (centered as  $X_0$  with 0 mean & rescaled with  $\text{Var}=1$ ) and  $Y$  (centered & rescaled as  $Y_0$ ):



Partial least squares (PLS) regression formulation:

$$\begin{cases} X_0 = X_S X_L^T + X_{\text{residuals}} \\ Y_0 = Y_S Y_L^T + Y_{\text{residuals}} \end{cases}$$

# Regression Modeling: PLSR

$$\begin{cases} X_0 = X_S X_L^T + X_{\text{residuals}} \\ Y_0 = Y_S Y_L^T + Y_{\text{residuals}} \end{cases}$$

observations (temporal)

$$\begin{matrix} & & \xrightarrow{\text{predictor variables (spatial)}} & & \\ \text{1st obs.} & \left[ \begin{array}{c|c|c|c} \text{1st predictor} & \cdots & & \text{p-th predictor} \\ \hline & X_0(1,:) & & \\ \vdots & \vdots & & \\ X_0(:,1) & \cdots & X \text{ (centered as } X_0) & \cdots & X_0(:,p) \\ \vdots & \vdots & & \vdots & \\ \hline & X_0(n,:) & & \end{array} \right] & = & \\ \vdots & & & & & & \\ \text{n-th obs.} & & & & & & \end{matrix}$$

observations (temporal)

$$\begin{matrix} & & \xrightarrow{\text{PLS components/modes number}} & & \\ \text{1st obs.} & \left[ \begin{array}{c|c|c|c} \text{1st PLS comp.} & \cdots & & \text{n}_{\text{comp}}\text{-th PLS comp.} \\ \hline & X_S(1,:) & & \\ \vdots & \vdots & & \\ X_S(:,1) & \cdots & X_S \text{ (X scores)} & \cdots & X_S(:,n_{\text{comp}}) \\ \vdots & \vdots & & \vdots & \\ \hline & X_S(n,:) & & \end{array} \right] & & \\ \vdots & & & & & & \\ \text{n-th obs.} & & & & & & \end{matrix}$$

$$\begin{matrix} & & \xrightarrow{\text{predictor variables (spatial)}} & & \\ \text{1st PLS comp.} & \left[ \begin{array}{c|c|c|c} \text{1st predictor} & \cdots & & \text{p-th predictor} \\ \hline & X_L^T(1,:) & & \\ \vdots & \vdots & & \\ X_L^T(:,1) & \cdots & X_L^T \text{ (X loading)} & \cdots & X_L^T(:,p) \\ \vdots & \vdots & & \vdots & \\ \hline & X_L^T(n_{\text{comp}},:) & & \end{array} \right] & & \\ \vdots & & & & & & \\ \text{n}_{\text{comp}}\text{-th comp.} & & & & & & \end{matrix}$$

response variables (spatial)

observations (temporal)

$$\begin{matrix} & & \xrightarrow{\text{response variables (spatial)}} & & \\ \text{1st obs.} & \left[ \begin{array}{c|c|c|c} \text{CO} & \cdots & & \text{NO}_2 \\ \hline & Y_0(1,:) & & \\ \vdots & \vdots & & \\ Y_0(:,1) & \cdots & Y \text{ (centered as } Y_0) & \cdots & Y_0(:,p) \\ \vdots & \vdots & & \vdots & \\ \hline & Y_0(n,:) & & \end{array} \right] & = & \\ \vdots & & & & & & \\ \text{n-th obs.} & & & & & & \end{matrix}$$

PLS components/modes number

response variables (spatial)

observations (temporal)

$$\begin{matrix} & & \xrightarrow{\text{PLS components/modes number}} & & \\ \text{1st obs.} & \left[ \begin{array}{c|c|c|c} \text{1st PLS comp.} & \cdots & & \text{n}_{\text{comp}}\text{-th PLS comp.} \\ \hline & Y_S(1,:) & & \\ \vdots & \vdots & & \\ Y_S(:,1) & \cdots & Y_S \text{ (Y scores)} & \cdots & Y_S(:,n_{\text{comp}}) \\ \vdots & \vdots & & \vdots & \\ \hline & Y_S(n,:) & & \end{array} \right] & & \\ \vdots & & & & & & \\ \text{n-th obs.} & & & & & & \end{matrix}$$

$$\begin{matrix} & & \xrightarrow{\text{response variables (spatial)}} & & \\ \text{1st PLS comp.} & \left[ \begin{array}{c|c|c|c} \text{CO} & \cdots & & \text{NO}_2 \\ \hline & Y_L^T(1,:) & & \\ \vdots & \vdots & & \\ Y_L^T(:,1) & \cdots & Y_L^T \text{ (Y loading)} & \cdots & Y_L^T(:,m) \\ \vdots & \vdots & & \vdots & \\ \hline & Y_L^T(n_{\text{comp}},:) & & \end{array} \right] & & \\ \vdots & & & & & & \\ \text{n}_{\text{comp}}\text{-th comp.} & & & & & & \end{matrix}$$

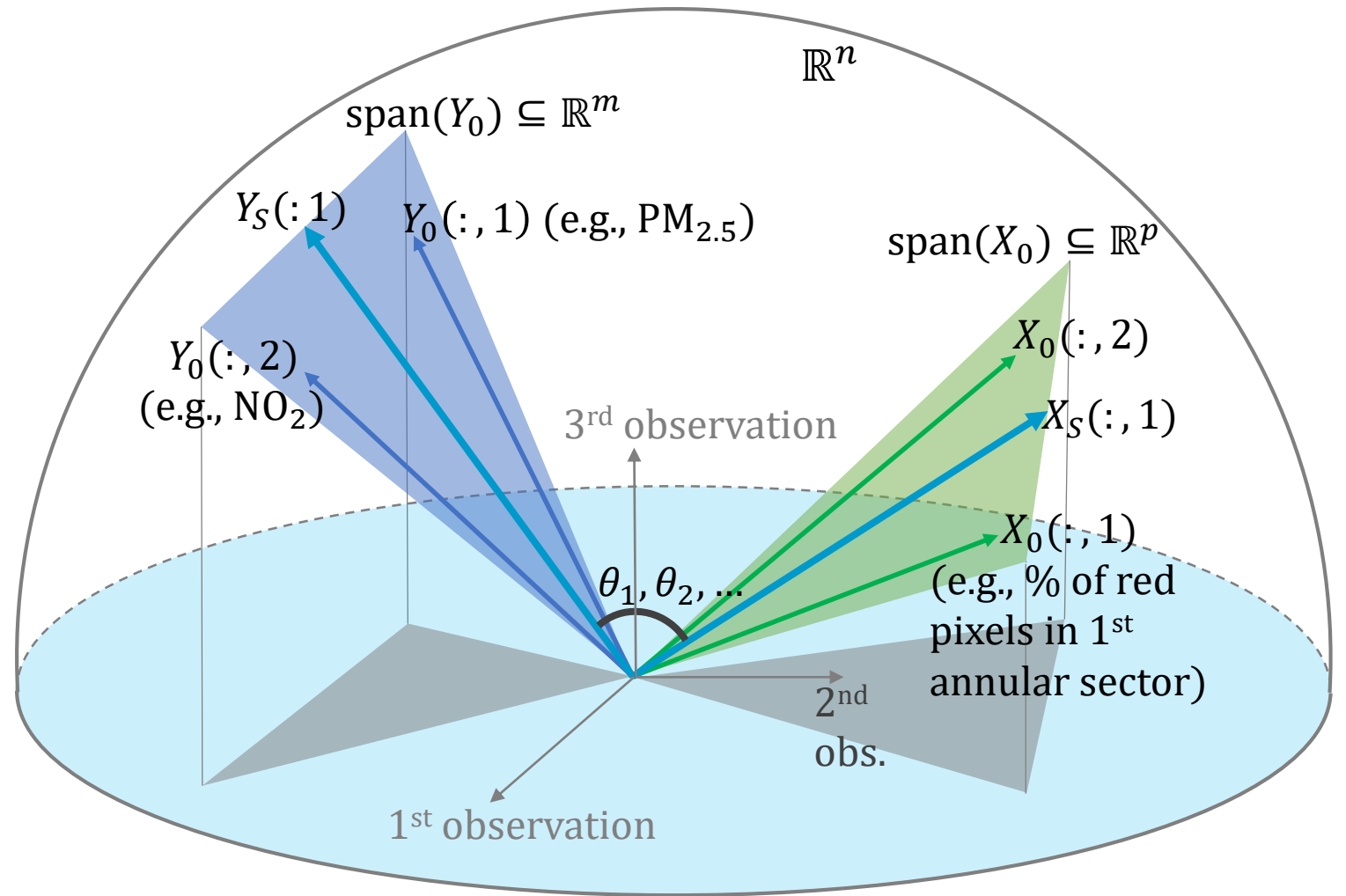
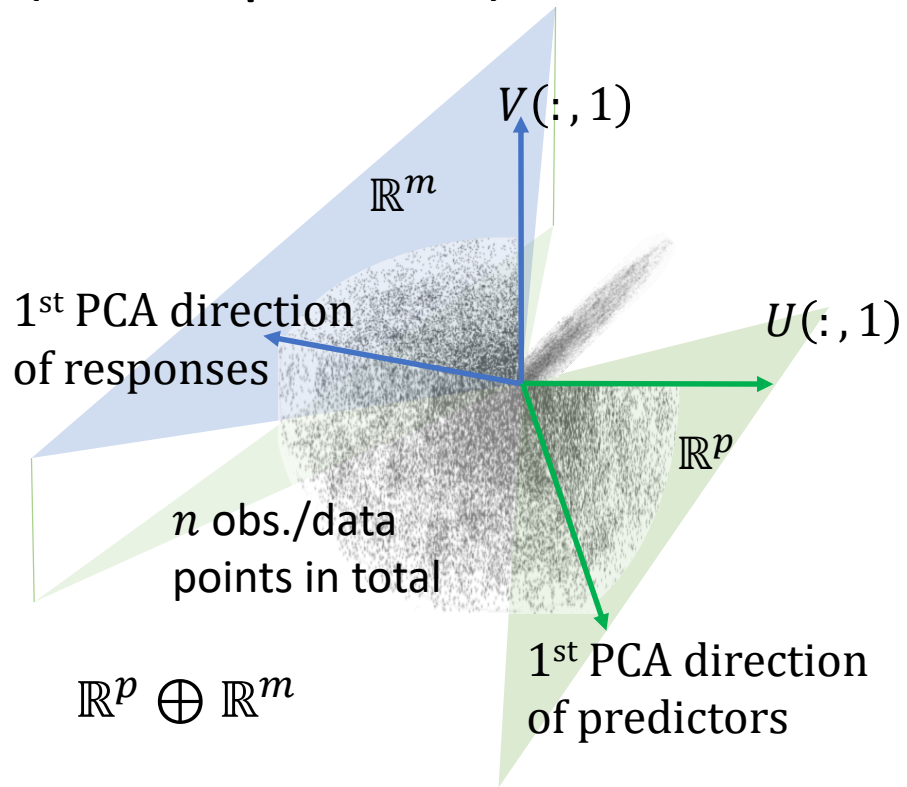
# Partial Least Squares (PLS)

- Partial least squares (PLS) regression formulation:

$$\begin{cases} X_0 = X_S X_L^T + X_{\text{residuals}} \\ Y_0 = Y_S Y_L^T + Y_{\text{residuals}} \end{cases}$$

- It aims to find correlation between  $X_0$  and  $Y_0$
- Recall SVD and covariance matrix ( $p$ -by- $m$ ):  $X_0^T Y_0 = U \Sigma V^T$
- Again, columns of  $U$  and  $V$  are orthonormal basis:  $U^T X_0^T Y_0 V = X_S^T Y_S = \Sigma$
- Projection of  $X_0$  on columns of  $U$  obtains predictor scores  $X_S$ ;
- Projection of  $Y_0$  on columns of  $V$  obtains response scores  $Y_S$ ;
- $X_S^T Y_S$  is inner products between columns of  $X_S$  and  $Y_S$
- If  $X_0$  and  $Y_0$  are centered and rescaled/normalized, projected on to  $U$  and  $V$  to get  $X_S$  and  $Y_S$ , so  $X_S^T Y_S = \Sigma$  contains cosines between columns of  $X_S$  and  $Y_S$ .
- $\Sigma$  is diagonal, non-increasing order: cosine (and hence correlation) between  $X_S(:, 1)$  and  $Y_S(:, 1)$  is maximal; angle is minimal

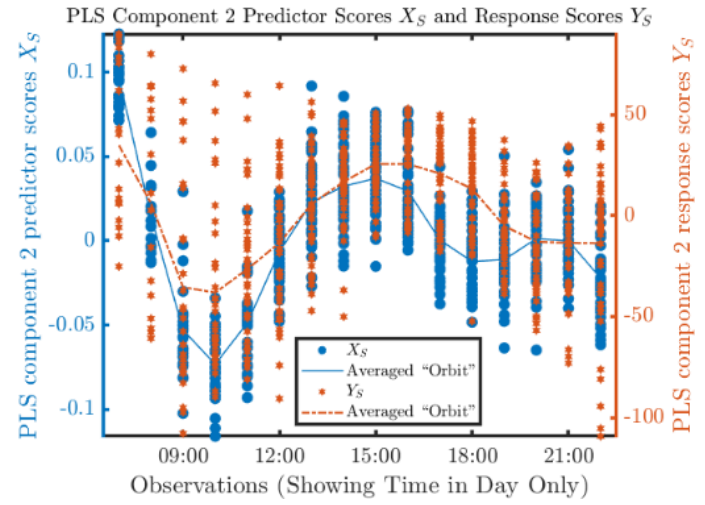
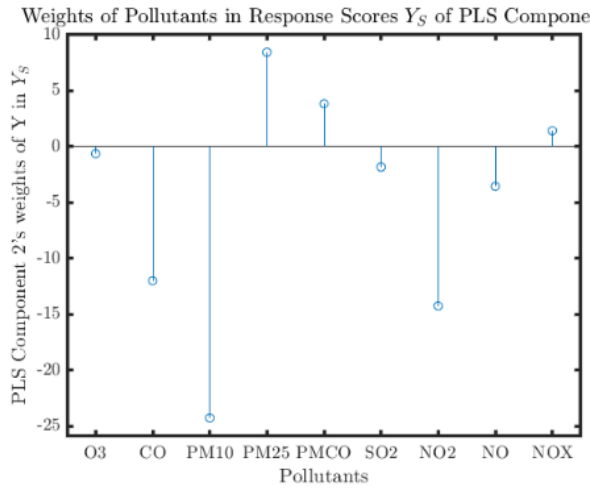
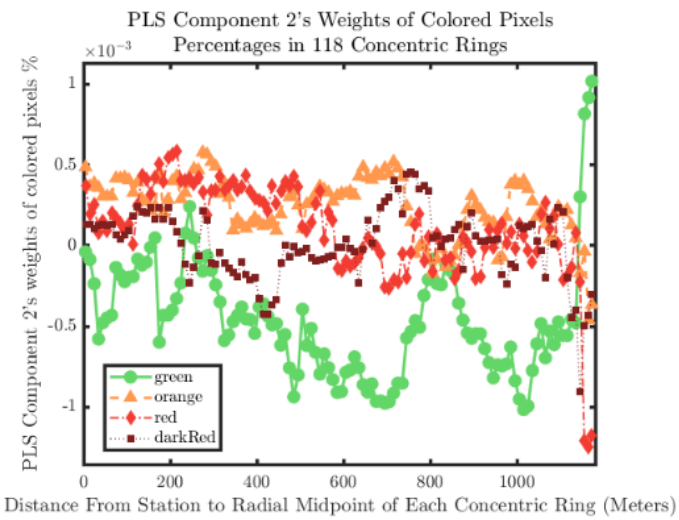
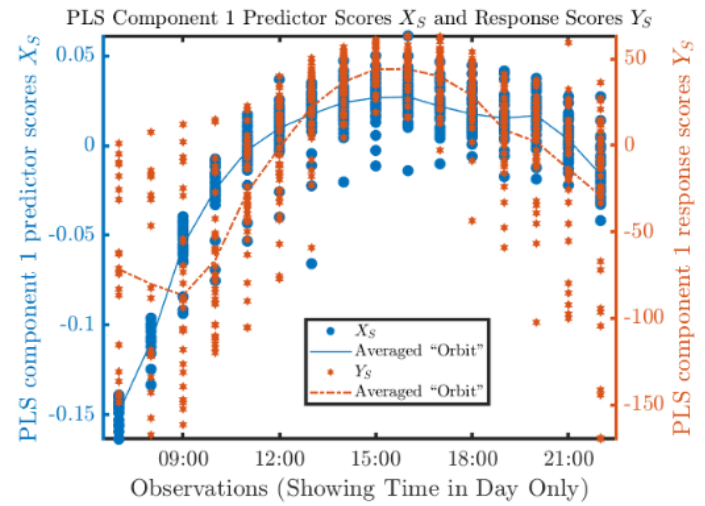
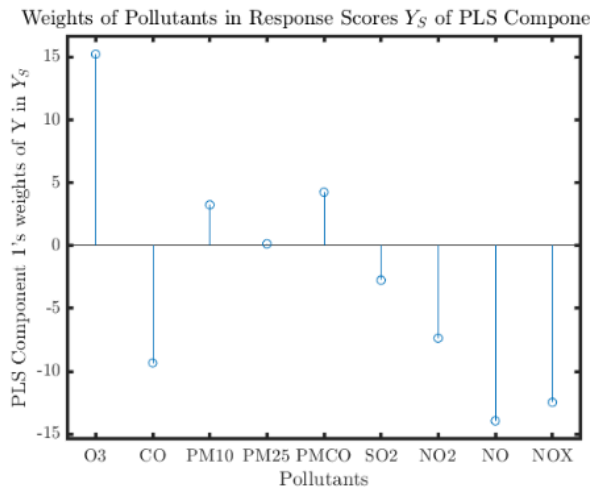
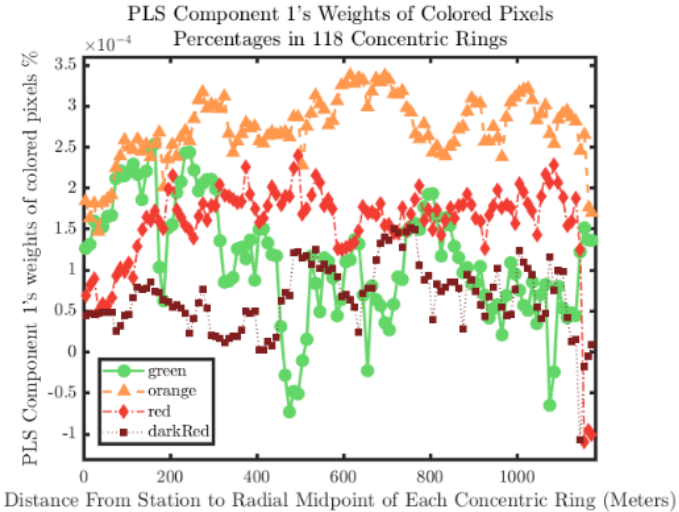
# Geometric Interpretation: principal angles between flats (subspaces)



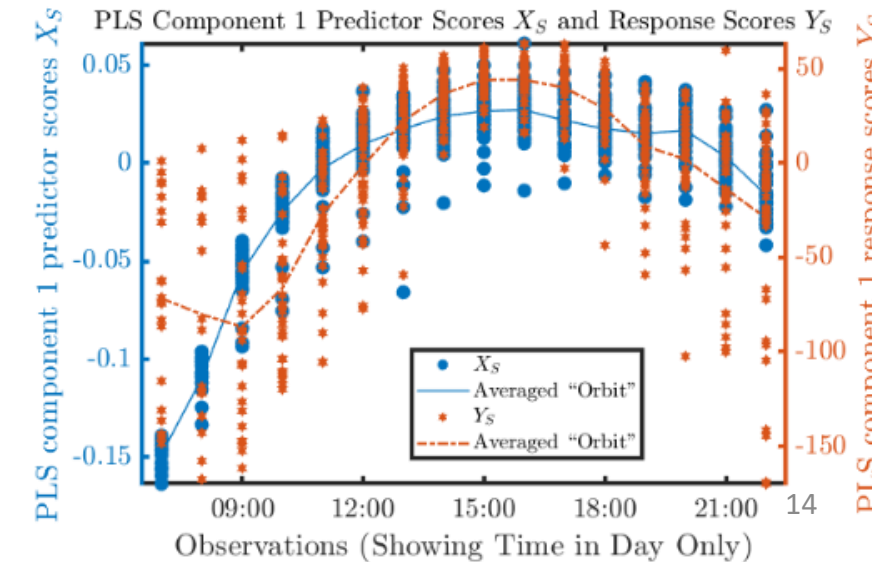
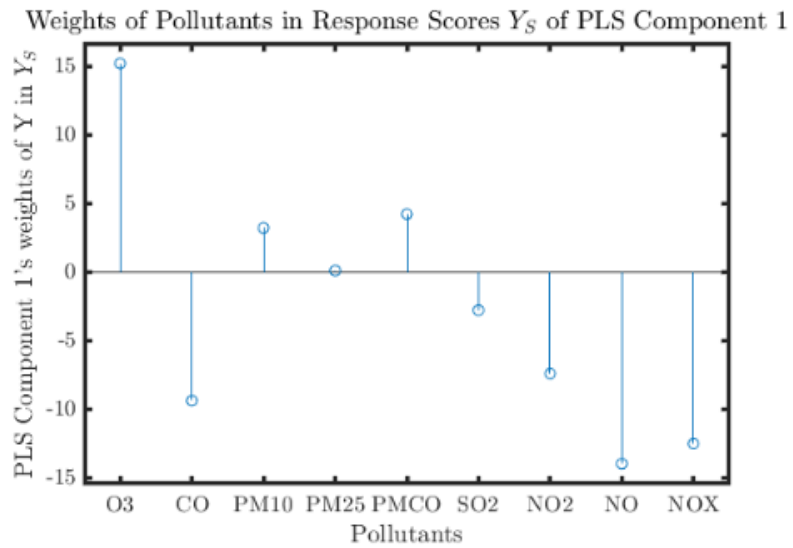
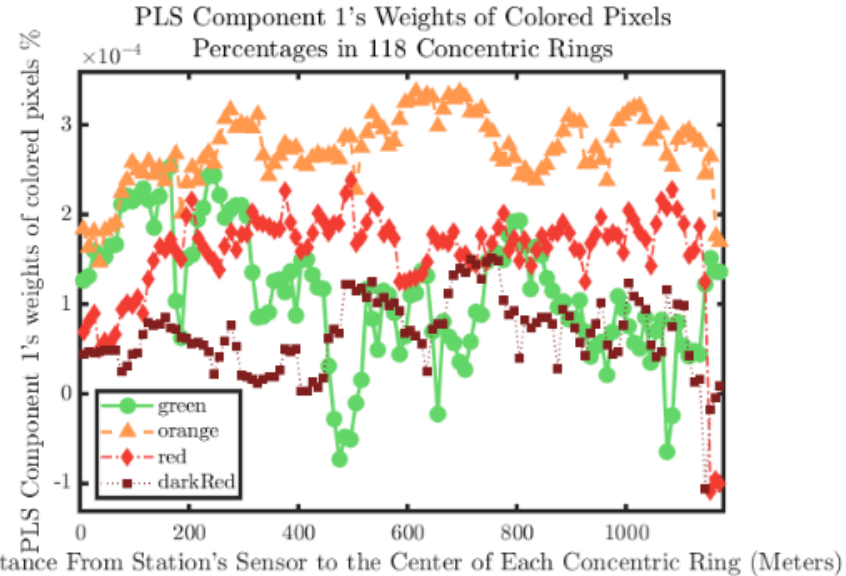
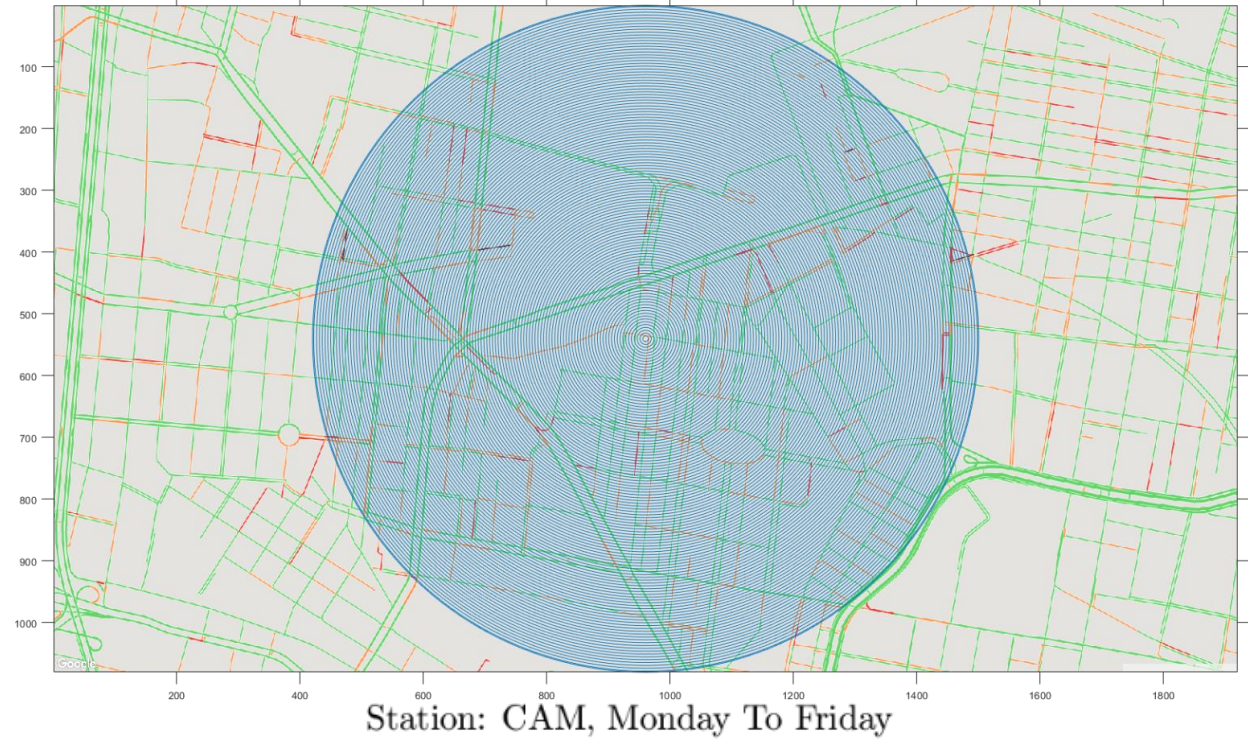
# Interpretations and Insight from PLS Regression Modeling

1st PLS component has physically meaningful interpretation:

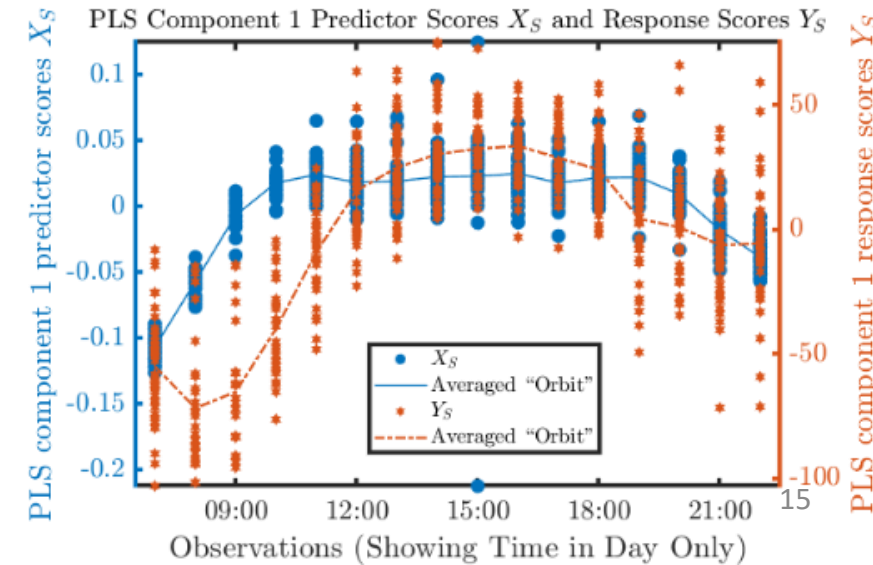
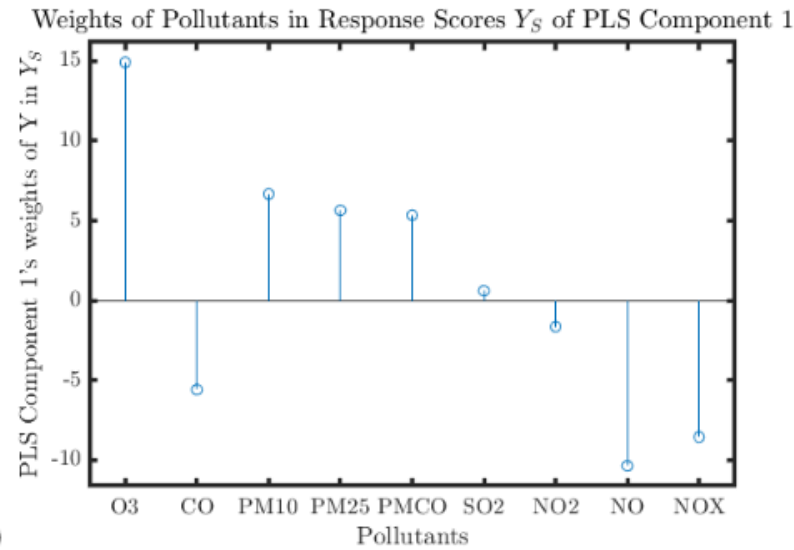
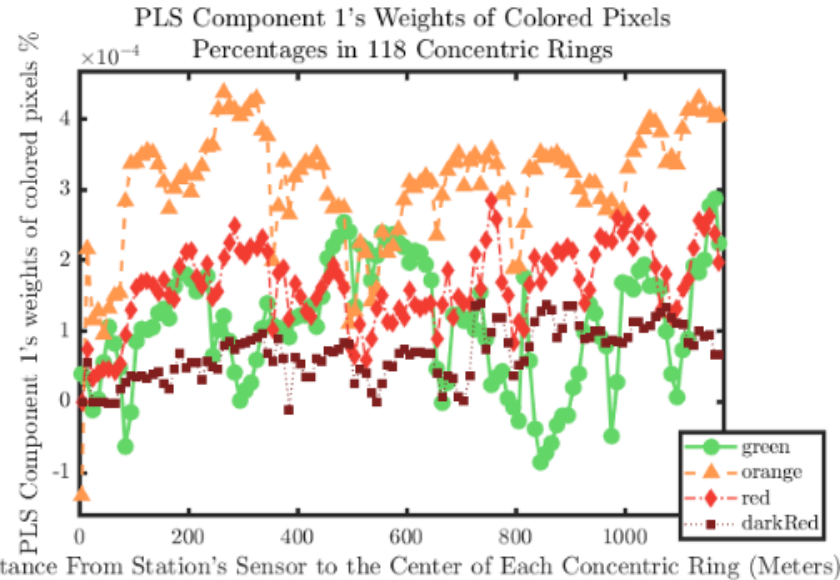
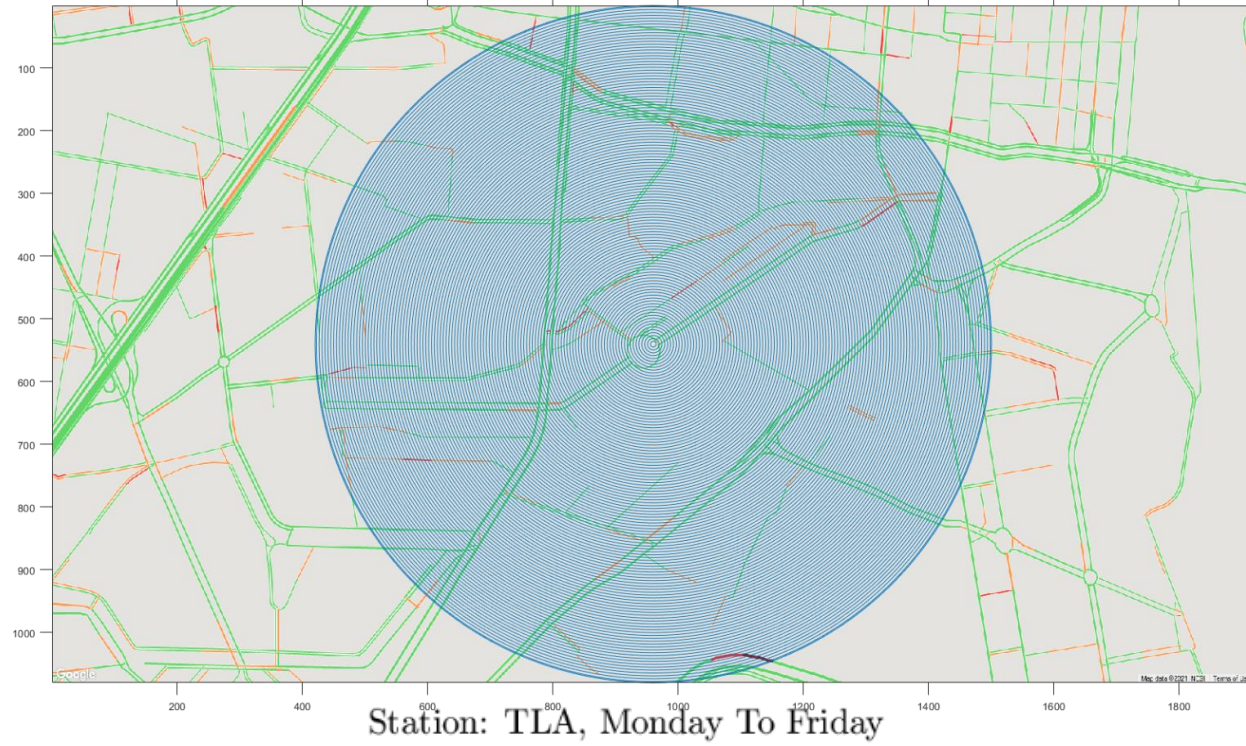
Station: CAM, Monday To Friday



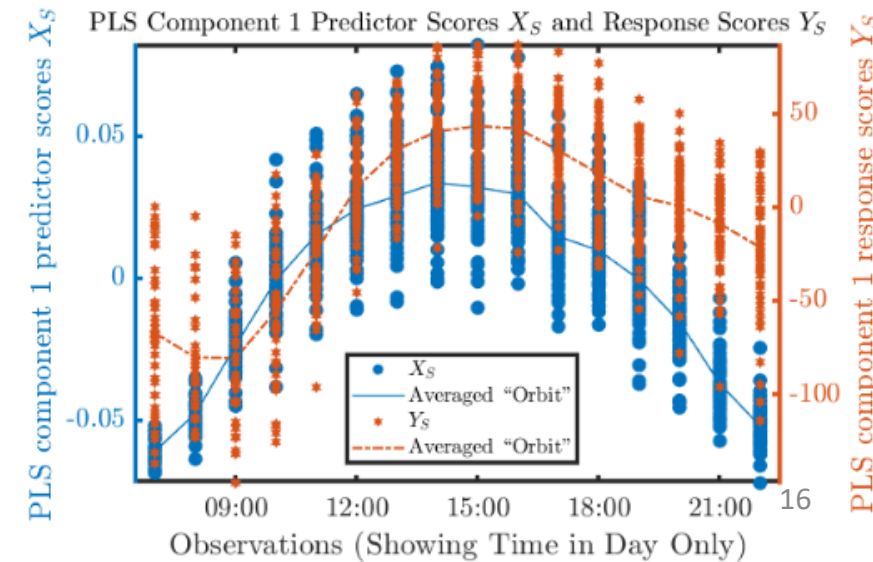
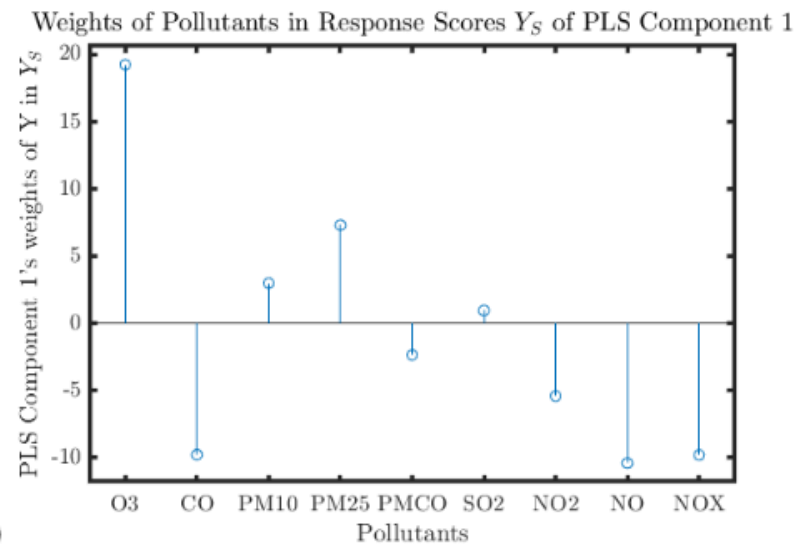
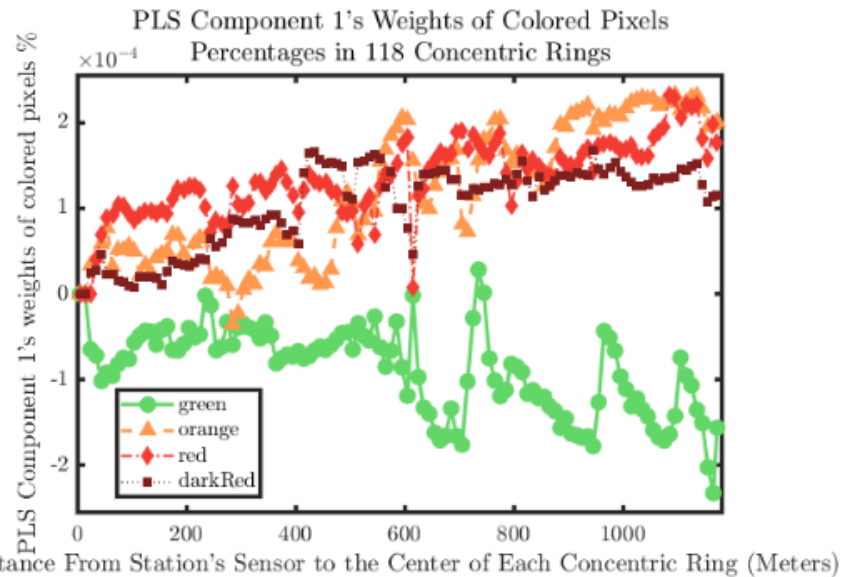
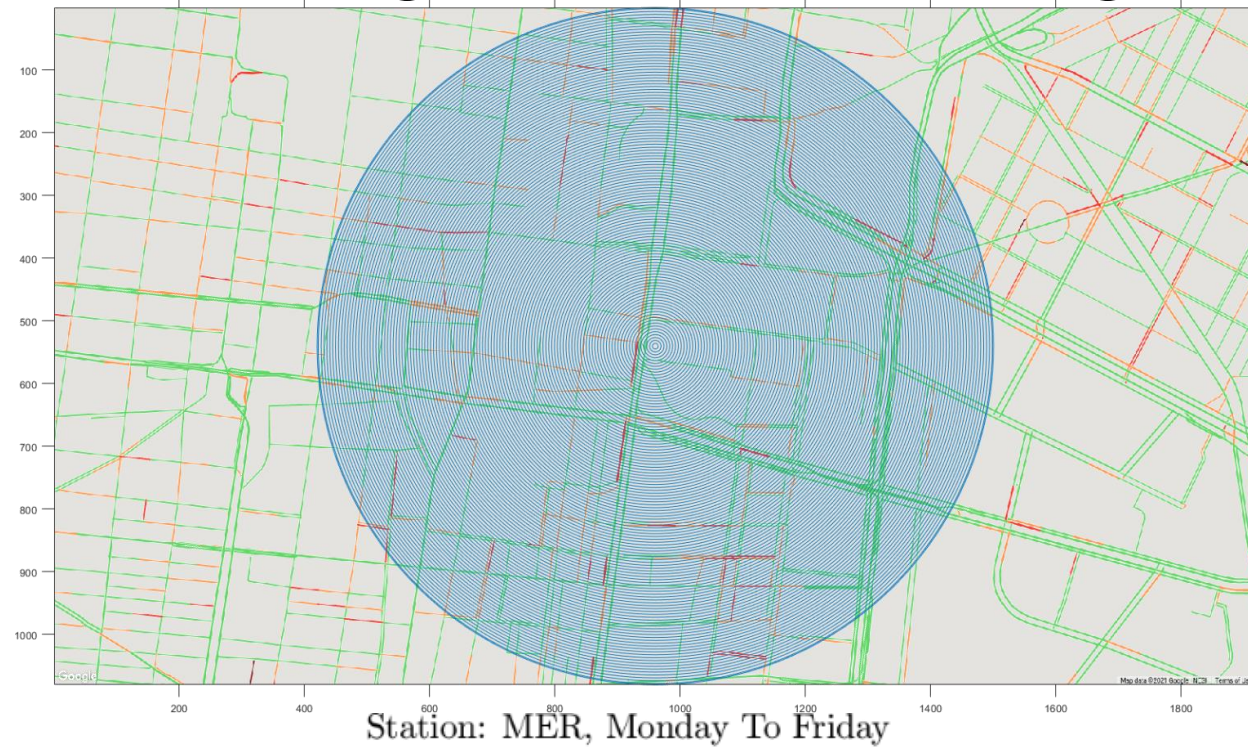
# Interpretations and Insight from PLS Regression Modeling



# Interpretations and Insight from PLS Regression Modeling



# Interpretations and Insight from PLS Regression Modeling





# Partial Least Squares (PLS) regression and PCA regression

- When # of response variables in  $Y$  is small, we can use PCA to reduce the spatial dimension of  $X$ , and then fit a least squares model (Principal component regression, PCR)
- When responses  $Y$  is multi-dimensional/high dimensional, PLS regression often performs better.
- In actual model fitting, a ( $p$  predictors)-by-( $m$  responses) coefficients matrix  $\beta_{n_{\text{comp}}}$  is fitted in least squares sense for

$$Y_S Y_L^T = X_S X_L^T \beta_{n_{\text{comp}}}$$

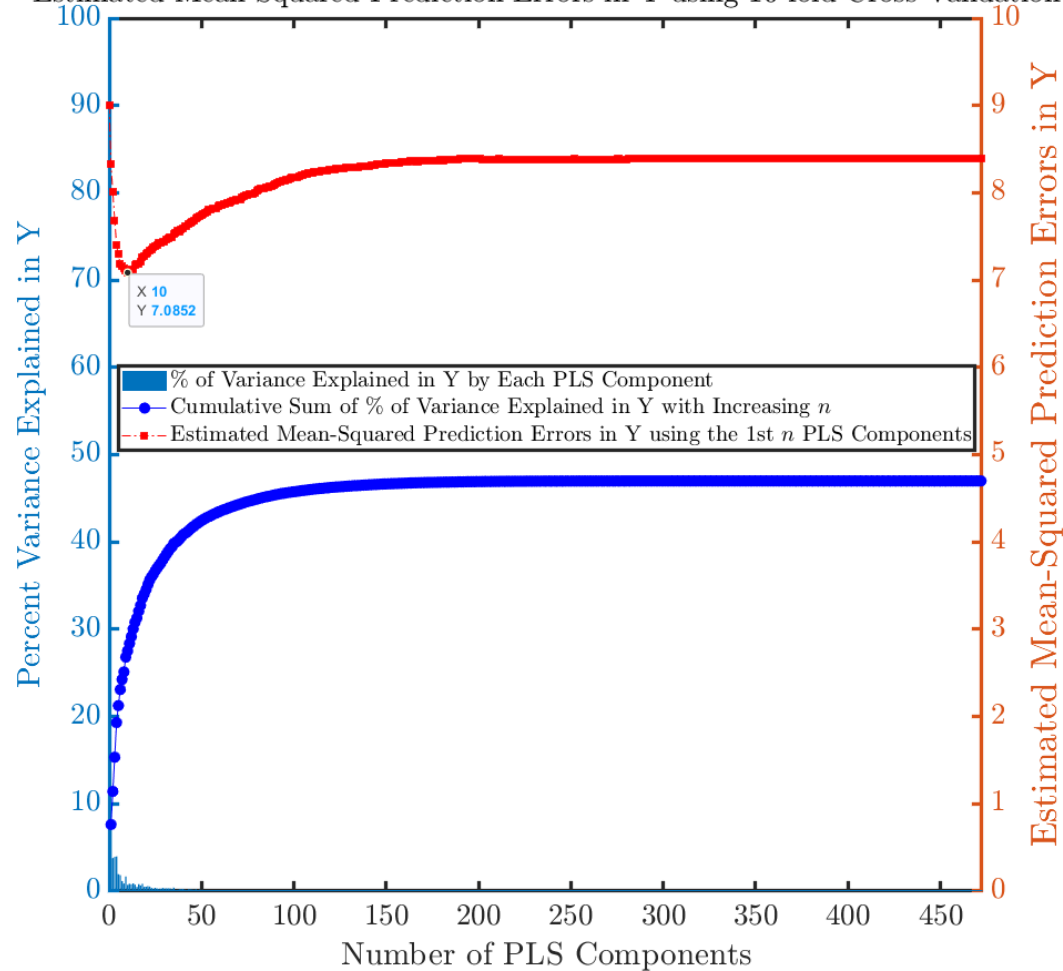
using truncated  $n_{\text{comp}}$  PLS components, and  $X_S X_L^T$  and  $Y_S Y_L^T$  are “reconstruction” of  $X_0$  and  $Y_0$

- Spatial dimension in  $X_0$  reduced from  $p$  to  $n_{\text{comp}}$ , effectively fitting  $Y_S = X_S \beta$  (a “partial” least-square), as compared to naively fit  $Y = X \beta$
- $n_{\text{comp}}$  can be fixed by cross-validation to minimize the expected mean-squared errors (MSE)  $\left(Y_0 - Y_S Y_L^T\right)^2$ .
- (a preliminary result on next page)

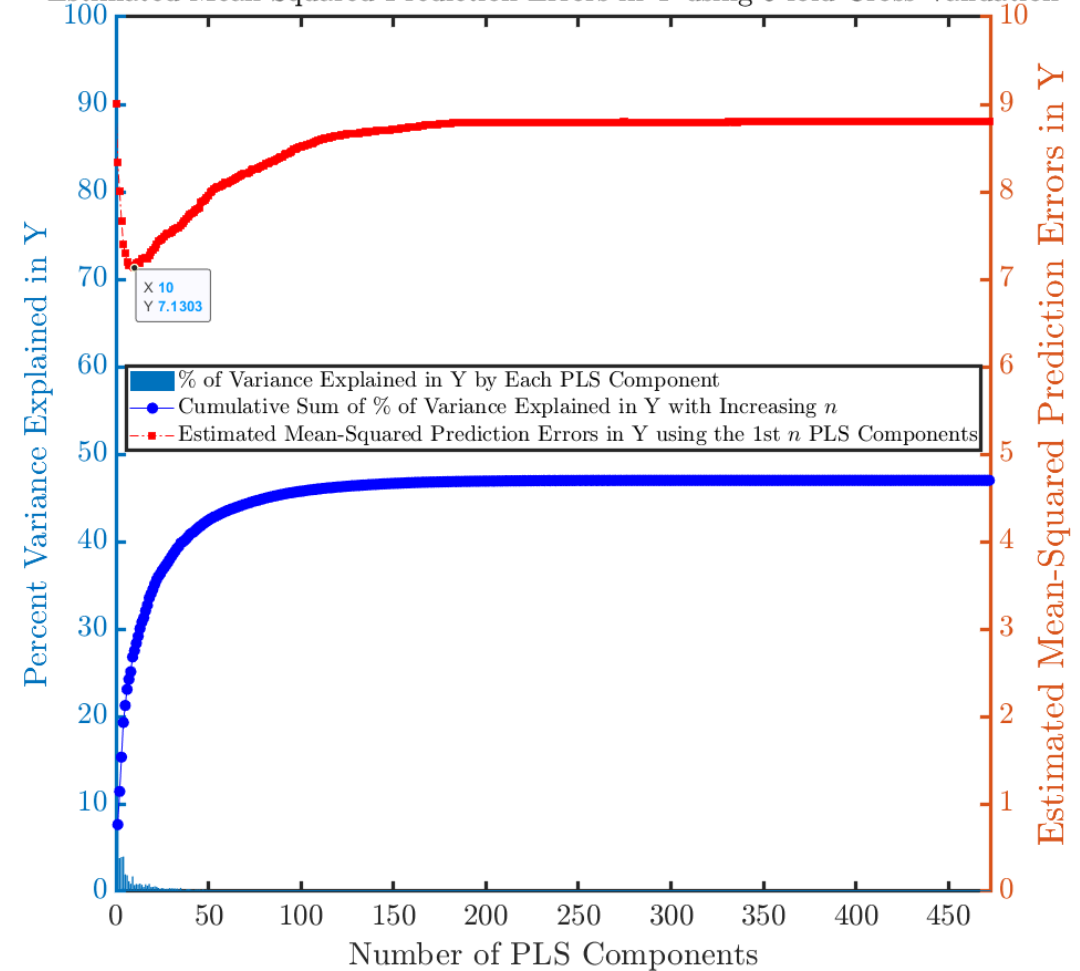
# PLSR Modeling and Prediction Performance

(all 9 pollutant response variables are centered to 0 mean with rescaled Var=1)

Station: 3 Stations Combined, Monday to Friday  
 $n = 472$  PLS components required to explain 47.0% of variance in Y  
 $\text{size}(X) = [2090, 472]$ ,  $\min(\text{size}(X,1)-1, \text{size}(X,2)) = 472$ ;  $\frac{472}{472} \approx 100.0\%$   
Estimated Mean Squared Prediction Errors in Y using 10-fold Cross-Validation



Station: 3 Stations Combined, Monday to Friday  
 $n = 472$  PLS components required to explain 47.0% of variance in Y  
 $\text{size}(X) = [2090, 472]$ ,  $\min(\text{size}(X,1)-1, \text{size}(X,2)) = 472$ ;  $\frac{472}{472} \approx 100.0\%$   
Estimated Mean Squared Prediction Errors in Y using 5-fold Cross-Validation



# Outlook and Conclusions

(conclusions by Marcella)

SAPIENS has built a database with both pollution measurements and traffic images, so we have:

- Cleaned and analysed the data and identified patterns
- Developed a model to extract the traffic intensities from Google Map images
- Used the regression modeling to (1) obtain interpretable insights on the relation between traffic and pollutants; and (2) train it on the data from three stations (traffic and pollution data) and cross-validated it to avoid overfitting

On-going activities:

- Validation/testing phase: use other sensors data to validate/test model
- Paper in preparation

# Outlook and Conclusions

There are more ideas and more possibilities to exploit and learn from these data.

More ideas on how to exploit the predicting power of the modeling

E.g., incorporating meteorological data, going beyond linear modeling techniques, etc.

Stay tuned for more from us

